# Machine Learning Week 8
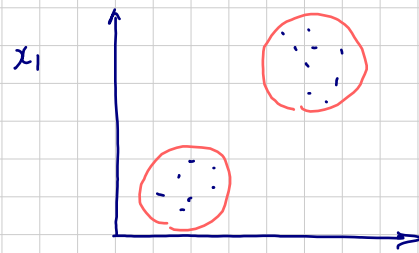
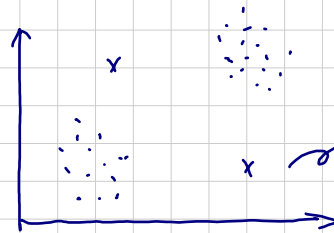## Unsupervised Learning

### 1. Clustering

$X_1$

group / clustering:
find some structures

Training set $\{x^{(1)}, x^{(2)}, \ldots, x^{(m)}\}$

$X_2$

Applications
- market segmentation
- social network analysis
- organize computer clusters
- astronomical data analysis

### 2. K-means Algorithm.

X

X

cluster centroids. (randomly initiated)

step (1)    Cluster assignment.
- binary assignment. of datasets. depending on proximity.

(2)    Move centroid.
- move to "mean" of location in all labelled

re-colour.

[Input].

① k ( number of clusters )

② Training set $\{ x^{(1)}, x^{(2)}, \ldots, x^{(m)} \}$

$$x^{(i)} \in \mathbb{R}^n \quad [\text{drop} \; x_0 = 1 \; \text{convention}]$$

[Algorithm]

Randomly initialize $K$ cluster centroids $\mu_k$, $k = 1 \cdots K$

do {

    for $i = 1 \cdots m$.

cluster assignment $\Big($
$$c^i := \min_k \| x^{(i)} - \mu_k \|^2$$
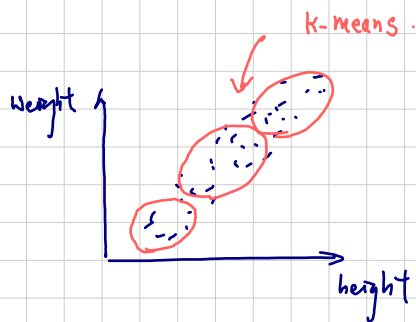
    for $k = 1 \cdots K$

move centroid $\Big($
$$\mu_k := \text{mean} ( \text{pts assigned to cluster } k )$$
$$\in \mathbb{R}^n$$

3. K-means for non-separated clusters.

weight

T-shirt sizing

K-means.



height

# Optimization Objective

K-means optimization objective.

- $c^i$ : index.

- $\mu_k$ : cluster centroid $k$     $\in \mathbb{R}^n$

- $\mu_{c^i}$ = cluster centroid of cluster to which $x^{(i)}$ has been assigned.

$$J(c^{(1)}, \ldots, c^{(m)}, \mu_1 \cdots \mu_k) = \frac{1}{m} \sum_{i=1}^{m} \| x^{(i)} - \mu_{c^{(i)}} \|^2$$

$$\min_{\substack{c^{(1)} \ldots c^{(m)} \\ \mu_1 \cdots \mu_k}} J(\cdots)$$

$x^{(i)}$

$\to x \quad \mu$.

min.

[Algorithm]

Randomly initialize K cluster centroids $\mu_k$. $k = 1 \cdots K$

do {

   for $i = 1 \cdots m$.

cluster assignment (

   $c^i := \min_k \| x^{(i)} - \mu_k \|^2$

   # Minimize $J(\cdots)$
   wrt.      $c^{(1)} \cdots c^{(m)}$.
   [holding $\mu_1 \cdots \mu_k$]

   for $k = 1 \cdots K$

move centroid (

   $\mu_k := \text{mean}(\text{pts assigned to cluster } k)$

   $\in \mathbb{R}^n$

   # Min. w.r.t. $\mu_1 \cdots \mu_k$

# Random Initialization

Rules:
    (1)   $K < m$

    (2)   Randomly pick $K$ training examples

    (3)   set $\mu_k$ = examples

· Might have different clustering (local optimum) $\longrightarrow$ try different random initialization.

- Implementation:

    for $i = 1 \cdots 1000$ {

        randomly initialize $k$-means
        run $k$-means. get $c^{(1)} \cdots c^{(m)}$, $\mu_1 \cdots \mu_k$
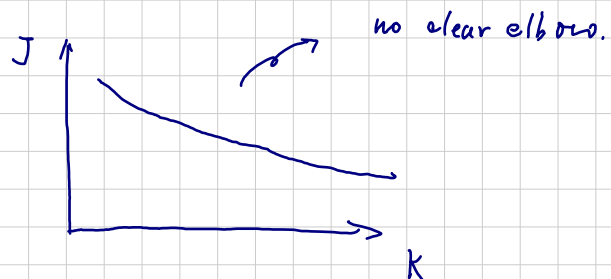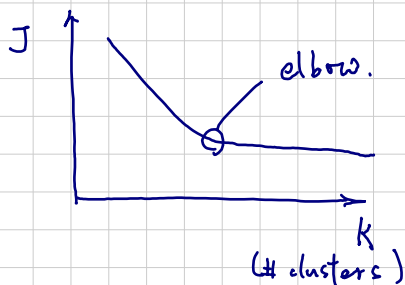
        complete cost function (distortion)
            $\hookrightarrow J(c^{(1)}, \cdots, c^{(m)}, \mu_1 \cdots \mu_k)$

    }

    Pick one with lowest cost $J$.
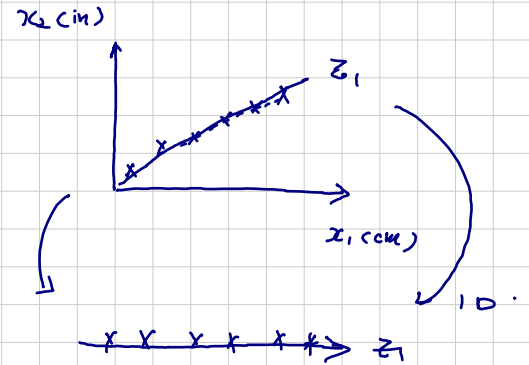
# Choose $K$

1.      Elbow method:



no clear elbow.

2.      Pre-defined ( t-shirt sizes : $s \cdots 5$ ).

# Dimensionality Reduction

1. Motivation I:   Data compression

   - Redundent data dimension ( cm, inch, ... )

$x_2$ (in)

$z_1$

$x_1$ (cm)

1 D.

$z_1$

ex. 3D → 2D.

$x^{(i)} \in \mathbb{R}^3$

$x_3$

[plane]

$x_2$

$x_1$

project onto a plane

(corelation approximation)

→ need 2-axis for plane

$z_1$

$z_2$

$z^{(i)} \in \mathbb{R}^2$.
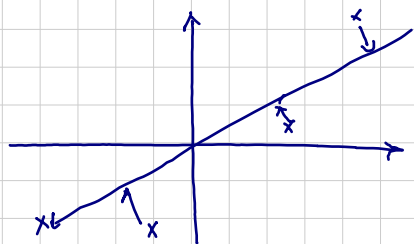
$\begin{pmatrix} z_1^{(i)} \\ z_2^{(i)} \end{pmatrix}$

procedure.    project data to one less dimension. ⟶ re-coordinate.

2  Motivation II:   Data visualization.

   . Combine certain features (corelated).    $x \to z$

                                              (500)    (2)

# Principal Component Analysis



Try to find a lower dimension surface to min. projection error.

⟶ perform mean normali. and feature scaling.

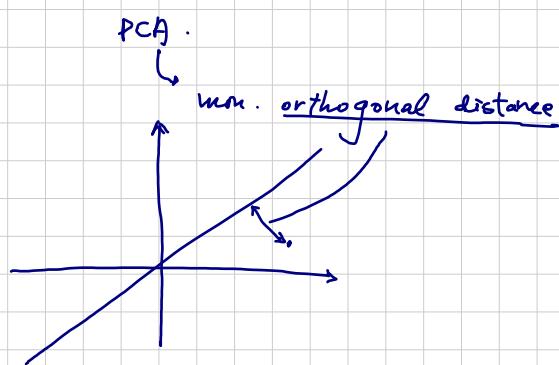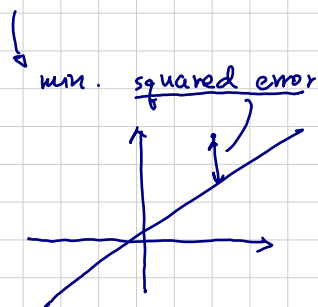. Principal Component Analysis [PCA] problem formulation

$2 \to 1$
- find dir. vector $u^{(1)} \in \mathbb{R}^k$ onto which data projects

$n \to k$

. find $k$ vectors $u^{(1)}, u^{(2)}, \ldots, u^{(k)}$ onto which data projects.

[project data to the linear subset span by $\{u^{(i)}\}_1^k$ ]

. linear regression

↳ min. <u>squared error</u>



PCA .
↳ min. <u>orthogonal distance</u>



Applys to higher dimension.

# Principal component analysis component.

· Data preprocessing

- mean normalisation
- feature scaling

· Reduction $\quad X^{(i)} \in \mathbb{R}^2 \longrightarrow z^{(i)} \in \mathbb{R} \qquad [\text{2D to 1D}]$

· Mathematical derivation (complicated)

## Algo.

① Reduce data from $\mathbb{R}^u$ to $\mathbb{R}^k$ $\qquad$ Sigma. $\in \mathbb{R}^{n \times n}$ $\qquad X = \left(\begin{array}{c}\frac{n}{=}\end{array}\right)$

$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad (n \times 1) \; (1 \times n)$

Compute "covariance matrix" $\qquad \Sigma = \frac{1}{m} \sum_{i=1}^{n} [x^{(i)}][x^{(i)}]^T = \left(\frac{1}{m}\right) \cdot X^T \cdot X.$

Compute "eigenvectors" of matrix $\Sigma$ $\qquad \longrightarrow$ singular value decomposition

$$[U, S, V] = svd(Sigma)$$

$$U = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \cdots & u^{(n)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

$\qquad \qquad \qquad$ first $k$ vectors : $u^{(1)} \cdots u^{(k)}$

② Projection.

$$U_{reduce} = \begin{bmatrix} | & & | \\ u^{(1)} & \cdots & u^{(k)} \\ | & & | \end{bmatrix}$$

$\qquad \qquad \qquad \qquad \in \mathbb{R}^{n \times k}$

$\qquad \qquad \qquad \qquad \qquad \qquad (k \times n) \qquad \qquad (n \times 1)$

$$z^{(i)} = U_{red}^T \cdot X^{(i)} = \left(\begin{array}{c} - \; u^{(1)} \; - \\ \vdots \\ - u^{(k)} \; - \end{array}\right) X^{(i)} \qquad = (k \times 1)$$

# Reconstruction from Compressed Representation

$$z = U_{red}^T x.$$

$$z \in \mathbb{R} \longrightarrow_{?} x \in \mathbb{R}^2$$

$$x \approx x_{approx} = \underbrace{U_{reduce}}_{(n \times k)} \cdot \underbrace{z}_{(k \times 1)}$$

## Choosing $k$ (# of principal component)

- Average Square projection error $\frac{1}{m} \sum_{i=1}^{m} \| x^{(i)} - x_{approx}^{(i)} \|^2$

- Total variation in data $\frac{1}{m} \sum_{i=1}^{m} \| x^{(i)} \|^2$

- Typically, choose $k$ to be min s.t.

$$\frac{\frac{1}{m} \sum_{i=1}^{m} \| x^{(i)} - x_{approx}^{(i)} \|^2}{\frac{1}{m} \sum_{i=1}^{m} \| x^{(i)} \|^2} \leq 0.01 \quad (1\%) \quad [*]$$

`` 99% variance retained ''

- Algorithm.

1) try PCA with $k = 1$

2) Compute $U_{reduce}, z^{(1)}, z^{(2)}, \dots z^{(m)}$
   $x_{approx}^{(1)} \dots x_{approx}^{(m)}$

3) check [*]

4) $k = k+1 \longrightarrow 2$

- $[U, S, V] = svd(sigma)$
  $\{$
  $n \times n$ diagonal, $S = \begin{pmatrix} S_{11} & & 0 \\ & \ddots & \\ 0 & & S_{nn} \end{pmatrix}$

For given $k$, (*) can be computed as $\left( 1 - \frac{\sum_{i=1}^{k} S_{ii}}{\sum_{i=1}^{n} S_{ii}} \right)$
$\underbrace{\qquad}_{\geq 0.99}$

# Advice for applying PCA

## Supervised Learning Speed Up

$$x^{(i)} \in \mathbb{R}^{10000}$$

$$(x^{(1)}, y^{(1)}) \ldots (x^{(m)}, y^{(m)})$$

Extract inputs:

Unlabelled dataset $x^{(1)}, x^{(2)} \ldots x^{(m)} \in \mathbb{R}^{10000}$

$\downarrow$ PCA   (defined using only on training set)

$$z^{(1)} \ldots z^{(m)} \in \mathbb{R}^{1000}$$

New training set

$$(z^{(1)}, y^{(1)}) \ldots (z^{(m)}, y^{(m)}) \longrightarrow h_\theta(z) = \frac{1}{1 + \exp(-\theta^T z)}$$

$$\text{train.}$$

Example

$$X \rightarrow z \longrightarrow h(z)$$
$$\text{(prediction)}$$

## Application

- Compression $\rightarrow$ choose k by % variance retain.
- Visualization $\rightarrow$ k = 2, 3

## What not to do $\longrightarrow$ prevent overfitting

. Use $z^{(i)}$ instead of $x^{(i)}$. reduce # of features from $n$ to $k$.

$\longrightarrow$ fewer features, less likely to overfit.   ( Why $\rightarrow$ PCA throws away information w/o knowing y )

. should use $\lambda$.

PCA usage.

Design of ML sys. ?

    -get training set

    [ _ run PCA reduce $x^{(i)} \rightarrow z^{(i)}$ ] do w/o this step first.

    _ train logistr.

    _ test set